RAG on GPU Max

| Sep_Llama2_GPU_PVC · · × + | | | | ~ | | × |
|---|----------|---|------|------------------|------|---|
| $\leftarrow \rightarrow C$ \bigcirc $\&$ \approx 70 | % ۲ | ያ | | ${igsidentials}$ | ப் | ≡ |
| 🔊 Rocky Linux 🔗 Rocky Wiki 🧖 Rocky Forums 🧖 Rocky Mattermost 🧔 Rocky Reddit | | | | | | |
| | | | i Ac | RUNNING . | Stop | : |
| | | | | | | |
| Lines Devices | | | | | | |
| 🦐 Llama Banker | | | | | | |
| Input your prompt here | <u>,</u> | | | | | |
| what was the FY2022 return on equity? | J | | | | | |
| | | | | | | |
| | | | | | | |
| | | | | | | |
| | | | | | | |
| | | | | | | |
| | | | | | | |
| | | | | | | |
| | | | | | | |
| | | | | | | |
| | | | | | | |
| | | | | | | |
| | | | | | | |

Introduction

In 2023, Intel launched its large-scale GPU, the Intel GPU Max, targeting high-performance computing (HPC) and Al applications at a cost-efficient price. Around the same time, OpenAl revolutionized the Al landscape with ChatGPT, reigniting the Al hype with generative AI (GenAl). As organizations began to adopt large language models (LLMs) like ChatGPT, they quickly recognized both the potential efficiency gains for employees and the risks associated with uncontrolled use. Publicly available LLMs can learn and retain information from users, posing a risk of inadvertently exposing confidential company data.

RAG to guardrail LLMs for enterprises

To address this concern, enterprises can deploy LLMs combined with Retrieval-Augmented Generation (RAG) on their own IT infrastructure, whether on-premises or in a cloud-based virtual data center. This approach ensures that the LLM is isolated from public use, safeguarding sensitive information. Additionally, RAG allows enterprises to provide secure, confidential information as a knowledge base to the LLM and establish guardrails without the need for expensive and resource-intensive model training or fine-tuning. Unlike training and fine-tuning, which require significant computational resources, RAG is less compute-intensive and more feasible for enterprises to implement.

Intel's GPU Max RAG

Since RAG is the lowest-hanging fruit to provide a safe and tailored LMMs use for enterprises, I built a RAG for Intel's GPU Max.

To do that , I adapted the code from Nicholas Renotte available here: <u>https://github.com/nicknochnack/Llama2RAG</u> That RAG implementation is intended for Nvidia GPUs, so I modified it to run on Intel GPU Max. You can see my code and tutorial here: <u>https://github.com/TheFavAI/Intel-GPU-Max-RAG</u>

A common misconception is that you need to convert CUDA code, running on Nvidia GPUs, into DPC++, running on Intel GPUs. The thing is that it isn't the case. Most AI Engineer, Data scientist, etc. code in Python. And by using Python, they usually don't need to code a single line of CUDA or DPC++. Indeed, everything is handled by the libraries used in their python codes!

The adaptions I had to do to the code of Nicolas Renotte are minimal.

Changes to the code

<u>Added line 5</u> : import intel_extension_for_pytorch as ipex

<u>Deletion of ", load_in_8bit=True" and addition of ".to("xpu")" line 25-27</u>: model = AutoModelForCausalLM.from_pretrained(name, cache_dir='./model/', use_auth_token=auth_token, torch_dtype=torch.float16, rope_scaling={"type": "dynamic", "factor": 2}).to("xpu")

<u>Added line 28</u>: model = ipex.optimize(model)

<u>Change to the pdf loader functions, but it has nothing to do with the RAG itself</u>: file_path_str = str(Path('./data/annualreport.pdf')) documents = loader.load(file_path=file_path_str, metadata=True)

How does it look?

Here is the code running. You can see the GPU being at 100% utilization to provide the answer as fast as possible.



The result

Here is an example of the answer the RAG can provider. The input is a company's financials, a huge pdf of more than 350 pages with a lot of numbers and complex financial notions. I asked an industry specific question to the RAG. And it manages to successfully answer as you can see.



Next steps

- As you can see, the UI isn't as good as it could be.
- The RAG only takes 1 document as an input, even though the document is huge. The RAG would need a database to handle multiple documents.